

AD-A198 406

Unclassified

DTIC FILE COPY

②

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

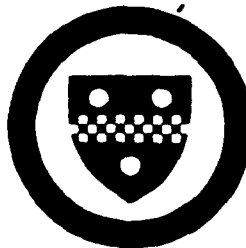
REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER <b>AFOSR-TK- 88-0796</b>	2. GOVT ACCESSION NO.	3. RECIPIENT'S CATALOG NUMBER
4. TITLE (and Subtitle) Detection of Change Points Using Rank Methods		5. TYPE OF REPORT & PERIOD COVERED <i>Report</i> Technical - March 1988
7. AUTHOR(s) B. Q. Miao and L. C. Zhao		6. PERFORMING ORG. REPORT NUMBER 88-02
9. PERFORMING ORGANIZATION NAME AND ADDRESS Center for Multivariate Analysis 515 Thackeray Hall University of Pittsburgh, PA 15260		8. CONTRACT OR GRANT NUMBER(s) AROS-88-0030
11. CONTROLLING OFFICE NAME AND ADDRESS Air Force Office of Scientific Research/ISRA 410 Dept. of the Air Force Bolling Air Force Base, DC 20332		10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS <i>6-1102P 7304 AD</i>
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office) <i>same as 11.</i>		12. REPORT DATE March 1988
		13. NUMBER OF PAGES 13
		15. SECURITY CLASS. (of this report) Unclassified
		16a. DECLASSIFICATION/DOWNGRADING SCHEDULE
16. DISTRIBUTION STATEMENT (of this Report)  Approved for public release; distribution unlimited.		
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)  <div style="text-align: right;">DTIC EXCE AUG 26 1988 D</div>		
18. SUPPLEMENTARY NOTES		
19. KEY WORDS (Continue on reverse side if necessary and identify by block number)  Change point, Detection, Directional data, Nonparametric method, rank statistics.		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number)  In this paper, the detection and estimation of change points of local parameters are studied by means of localization procedures and rank statistics. These techniques are also applied to detection and estimation of the change points of scale parameters and that of location parameters of directional data.		

DD FORM 1 JAN 73 1473

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

88 8 25 146

**Center for Multivariate Analysis**  
**University of Pittsburgh**



# DETECTION OF CHANGE POINTS USING RANK METHODS

B. Q. Miao and L. C. Zhao

Center for Multivariate Analysis  
University of Pittsburgh

Technical Report No. 88-02

February 1988

Center for Multivariate Analysis  
Fifth Floor Thackeray Hall  
University of Pittsburgh  
Pittsburgh, PA 15260

Accession For	
NTIS GRA&I	<input checked="checked" type="checkbox"/>
DTIC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	
By	
Distribution/	
Availability Codes	
Dist	Avail and/or Special
A-1	

Research Sponsored by the Air Force Office of Scientific Research under Grant AF60 - 88 - 0030. The United States Government is authorized to reproduce and distribute reprints for governmental purposes notwithstanding any copyright notation hereon.

# DETECTION OF CHANGE POINTS USING RANK METHODS

B.Q. Miao and L. C. Zhao

Center for Multivariate Analysis  
University of Pittsburgh  
Pittsburgh PA 15260

AMS 1980 Subject Classifications: Primary 62G30, Secondary 62D05.

**Key words and phrases:** Change point, Detection, Directional data, Nonparametric method, rank statistics.

## ABSTRACT

In this paper, the detection and estimation of change points of local parameters are studied by means of localization procedures and rank statistics. These techniques are also applied to detection and estimation of the change points of scale parameters and that of location parameters of directional data.

## 1. INTRODUCTION

Change point problem arises in many fields and attracts the attention of many authors. The techniques employed to detect and estimate the change point can generally be classified into two categories : parametric ( Krishnaiah and Miao .1988 ), and

nonparametric (Csörgö and Horváth, 1988). Bayesian methods also plays a major role (Zacks, 1983).

In this paper, we concentrate our attention on the problem of detection of change points of location parameter by localization and rank statistics when data are large. Our method is different from, and has some advantages over, the existing methods, such as CUSUM (cumulative sum) and Csörgö and Horváth's non-sequential nonparametric AMOC (at most one change) procedures. First, localized procedures reduce computation. Second, these detecting and estimating procedures require no moment condition, instead we only assume that observed data come from a continuous distribution with a unique median.

## 2. MODEL AND DETECTING AND ESTIMATING PROCEDURE

Let  $x(t)$  be an independent process on the interval  $(0,1]$  whose marginal distributions differ only by location parameters. Specifically speaking, there exist  $t_0, t_1, \dots, t_{q+1}$  and  $a_1, \dots, a_{q+1}$  such that  $0 = t_0 < t_1 < \dots < t_q < t_{q+1} = 1$  and

$$x(t) \sim F(x-a_j), \quad \text{if } t_{j-1} < t \leq t_j, j = 1, 2, \dots, q+1, \quad (2.1)$$

with

$$a_j \neq a_{j+1} \quad \text{for } j = 1, \dots, q. \quad (2.2)$$

As usual,  $t_1, \dots, t_q$  are called change points. In practical applications, the number  $q$  and the locations of change points are unknown. To estimate  $q$  and  $t_1, \dots, t_q$ , we sample this process sequentially at equal distances, and get  $x(1/n), x(2/n), \dots, x(n/n)$ . By assumption,  $x(1/n), \dots, x(n/n)$  are independent. Define  $k_0^\circ = 0$ ,  $k_{q+1}^\circ = n$ , and  $k_1^\circ, \dots, k_q^\circ$  as follows:

$$|k_i^\circ(n)/n - t_i| < 1/n, \quad i = 1, \dots, q. \quad (2.3)$$

Then we have

$$x(i/n) \sim F(x-a_j), \quad \text{if } k_{j-1}^\circ(n) < i \leq k_j^\circ(n), \quad j = 1, 2, \dots, q+1, \quad (2.4)$$

From this fact, to estimate  $(t_1, \dots, t_q)$  is equivalent to estimating  $(k_1^\circ(n)/n, \dots, k_q^\circ(n)/n)$ . For simplicity, set  $x(i/n) = x_i$  and  $k_i^\circ(n) = k_i^\circ$ . Hereafter,  $\alpha_n \gg \beta_n$  means  $\alpha_n/\beta_n \rightarrow \infty$ . Let  $m = m_n$  and  $C_n$  be positive integers such that  $n \gg m_n \gg C_n \gg \log n$ .

Set  $A_{k,m} = \{x_{k-m+1}, \dots, x_k, x_{k+1}, \dots, x_{k+m}\}$ ,  $k = m, \dots, n-m$ ,  $k_j$  be the rank of  $x_j$  in  $A_{k,m}$ .  
Define

$$y_{k,j} = \begin{cases} 1 & \text{if } x_{k-m+j} \leq x_{k_j} \\ -1 & \text{otherwise} \end{cases}, j = 1, \dots, m. \quad (2.5)$$

$$S_{k,m} = \sum_{j=1}^m y_{k,j}. \quad (2.6)$$

$$D_n = \{k : k = m, \dots, n-m : S_{k,m}^2/m > C_n\}, \quad (2.7)$$

$$k_{1,n} = \min\{k : k \in D_n\},$$

$$D_{1,n} = \{k : k \in D_n, k - k_{1,n} < 3m\}.$$

Next, put

$$k_{2,n} = \min\{k : k \in D_n - D_{1,n}\}$$

$$D_{2,n} = \{k : k \in D_n - D_{1,n}, k - k_{2,n} < 3m\}.$$

Continuing this process, we can define  $D_{2,n}, D_{2,n}, \dots$ , which are easily seen to be nonempty. We have

$$D_n = D_{1,n} + D_{2,n} + \dots + D_{\hat{q},n}.$$

Define

$$\hat{t}_j = 2^{-1}[k_{j,n} + \max(k : k \in D_{j,n})], \quad j = 1, 2, \dots, \hat{q}.$$

We have the following theorem.

**THEOREM.** If the distribution  $F$  is continuous and has unique median, then  $(\hat{q}, \hat{t}_1, \dots, \hat{t}_{\hat{q}})$  is a strongly consistent estimate of  $(q, t_1, \dots, t_q)$ .

### 3. LEMMAS

To prove the above theorem, we need some lemmas.

**LEMMA 1.** (Hoeffding, 1963) Let the population  $C$  consists of  $N$  values  $C_1, \dots, C_N$ . Let  $x_1, \dots, x_n$  denote a random sample without replacement from  $C$  and let  $y_1, \dots, y_n$  denote a random sample with replacement from  $C$ . If the function  $f(x)$  is continuous and convex then

$$Ef\left(\sum_{i=1}^n x_i\right) \leq Ef\left(\sum_{i=1}^n y_i\right).$$

**LEMMA 2.** Let the notation be defined as in section 2 and define

$$B_{k,m} = \{x_{k-m+1}, \dots, x_k\}, \quad k = m, m+1, \dots, n. \quad (3.1)$$

$$S_{k,m} = \sum_{j=1}^m y_{k-m+j}, \quad (3.2)$$

$$B_n = \{k : \exists j, 1 \leq j \leq q+1, \text{ such that } k_{j-1}^0 + 1 \leq k - m + 1 < k \leq k_j^0\}. \quad (3.3)$$

Then, we have

$$S_{k,m}^2/m = O(\log n) \quad \text{a.s.}$$

uniformly for all  $k \in B_n$ .

*Proof.* Let  $z_{k,1}, \dots, z_{k,m}$  be a random sample with replacement from population  $\{1, \dots, 1, -1, \dots, -1\}$ , where the number of 1's and -1's are both  $m$ . By lemma 1 for any  $t \in (0, 1/4)$  and  $A > 0$ , we have

$$\begin{aligned} P\{S_{k,m} \geq A \sqrt{m \log n}\} &\leq \exp\{-tA \sqrt{m \log n}\} E e^{t S_{k,m}} \\ &\leq \exp\{-tA \sqrt{m \log n}\} E \exp\left\{t \sum_{j=1}^m z_{k,j}\right\} \end{aligned}$$

$$\begin{aligned}
&= \exp\{-tA\sqrt{m\log n}\} (E \exp\{tz_{k,j}\})^m \\
&= \exp\{-tA\sqrt{m\log n}\} ((e^t + e^{-t})/2)^m \\
&\leq \exp\{-tA\sqrt{m\log n} + mt^2 e^t/2\}.
\end{aligned}$$

Since  $m \gg \log n$ , it is possible for  $n$  large to take  $t = A\sqrt{(\log n)/m} < 1/4$ . It follows

$$P(S_{k,m} \geq A\sqrt{m\log n}) < \exp\{-0.3A^2 \log n\}. \quad (3.4)$$

A similar argument gives

$$P(S_{k,m} \leq -A\sqrt{m\log n}) < \exp\{-0.3A^2 \log n\}. \quad (3.5)$$

The inequalities (3.3) and (3.4) imply

$$\sum_{n=1}^{\infty} P(\sup_{k \in B_n} |S_{k,m}| \geq A\sqrt{m\log n}) \leq \sum_{n=1}^{\infty} 2n \exp\{0.3A^2 \log n\}.$$

This is a convergence series if  $0.3A^2 > 2$ , Take  $A = 3$ , by Borel-Cantelli lemma, with probability one for large  $n$ , we have uniformly for all  $k \in B_n$ ,

$$|S_{k,m}| \leq 3\sqrt{m\log n},$$

$$S_{k,m}^2/m \leq 9\log n \quad \text{a.s.}$$

uniformly for all  $k \in B_n$ .

**LEMMA 3.** Let  $x_{i,k}^*$  denote the  $i$ -th sample order statistic in  $B_{k,m}$ . Let  $[\alpha m]$  denote the integer part of  $\alpha m$ . Assume that the  $\alpha$ -quantile of the continuous distribution  $F$  is unique, and is denoted by  $\mu_\alpha$ . Then

$$x_{[\alpha m],k}^* \rightarrow \mu_\alpha + a_j, \quad \text{a.s.}$$

uniformly for all  $k \in B_n$ , where  $k_{j-1}^0 + 1 \leq k - m + 1 < k \leq k_j^0$ .



*Proof.* Without loss of generality, we assume that  $0 < \alpha < 1$ ,  $k_{j-1}^0 + 1 \leq k - m + 1 < k \leq k_j^0$  for some fixed  $j$  and  $a_j = 0$ . Write  $r = [\alpha m]$ , then for any  $\varepsilon > 0$ ,

$$\begin{aligned} P(|x_{[\alpha m], k}^* - \mu_\alpha| \geq \varepsilon) &= P(x_{r, k}^* \leq \mu_\alpha - \varepsilon) + P(x_{r, k}^* \geq \mu_\alpha + \varepsilon) \\ &\triangleq I_1 + I_2. \end{aligned}$$

Set  $F(\mu_\alpha - \varepsilon) = \alpha - \delta$ . By the uniqueness of the  $\alpha$ -quantile of  $F$ ,  $\delta > 0$ . Without loss of generality, we can assume that  $\alpha - \delta > 0$ . Since  $F(x)$  is continuous, one gets,

$$\begin{aligned} I_1 &= P(x_{r, k}^* \leq \mu_\alpha - \varepsilon) \\ &= \frac{m!}{(r-1)!(m-r)!} \int_0^{F(\mu_\alpha - \varepsilon)} t^{r-1}(1-t)^{m-r} dt \\ &= \frac{r \cdot m!}{r!(m-r)!} \int_0^{\alpha - \delta} t^{r-1}(1-t)^{m-r} dt \end{aligned}$$

For a fixed  $\theta \in (0, 1)$ , put

$$g_\theta(t) = t^\theta(1-t)^{1-\theta}.$$

It is easy to see that  $g(t)$  is increasing in  $t$  on the interval  $[0, \theta]$ , which implies that  $t^{r-1}(1-t)^{m-r}$  is increasing for  $t \in [0, (r-1)/(m-1)]$ . Since  $r = [\alpha m]$ , we get  $\alpha - \delta < (r-1)/(m-1)$  for large  $m$ . Using Stirling's formula, we obtain

$$\begin{aligned} I_1 &\leq \frac{rm!}{r!(m-r)!} (\alpha - \delta)^{r-1} (1 - \alpha + \delta)^{m-r} \\ &\leq \frac{C_1 m \sqrt{2\pi m}}{\sqrt{2\pi \cdot 2\pi \cdot r(m-r)}} \frac{(\alpha - \delta)^{r-1} (1 - \alpha + \delta)^{m-r}}{(r/m)^r ((m-r)/m)^{m-r}} \\ &\leq C \sqrt{m} q_1^m, \end{aligned} \tag{3.6}$$

where  $C_1$  and  $C$  are constants depending on  $\alpha, \delta$ , and

$$\begin{aligned} q_1 &= \{(\alpha - \delta)/(\alpha - \delta/2)\}^\alpha \{(1 - \alpha + \delta)/(1 - \alpha + \delta/2)\}^{1-\alpha}, \\ &= g_\alpha(\alpha - \delta)/g_\alpha(\alpha - \delta/2) < 1. \end{aligned}$$

Similarly, one can show that

$$P(x_{r,k}^* \geq \mu + \varepsilon) \leq C \sqrt{m} q_2^m, \quad (3.7)$$

where  $0 < q_2 < 1$ . Take  $q = \max(q_1, q_2)$ , then (3.6), (3.7) and  $m_n \gg \log n$  together imply

$$\sum_{n=1}^{\infty} P(\sup_{k \in B_n} |x_{[\alpha m],k}^* - \mu_\alpha - a_j| \geq \varepsilon) \leq \sum_{n=1}^{\infty} 2cn\sqrt{m} q^m < \infty.$$

By Borel-Cantelli lemma, it follows that

$$x_{[\alpha m],k}^* \rightarrow \mu_\alpha + a_j, \text{ a.s.}$$

uniformly for all  $k \in B_n$ , where  $k_{j-1}^0 + 1 \leq k - m + 1 < k \leq k_j^0$ .

**LEMMA 4.** Let  $x_i \in A_{k_j^0, m}$ ,  $j = 1, 2, \dots, q$ ,  $S_{k_j^0, m}$  be as defined by (2.6). Then there exists a positive  $\lambda$  such that with probability one for  $n$  large,

$$S_{k_j^0, m}^2 / m \geq m\lambda^2, \quad j = 1, \dots, q.$$

*Proof.* Let  $x_i \in A_{k_j^0, m}$  for some fixed  $j$ . Without loss of generality, we assume  $a_{j+1} > a_j$ . Take  $\lambda \in (0, 1/2)$  satisfying the following conditions:

1° the  $(1/2 - \lambda)$  and  $(1/2 + \lambda)$  - quantiles of the continuous distribution  $F(x)$  are unique.

$$2^\circ \quad \mu_{1/2} - \mu_{1/2-\lambda} < (a_{j+1} - a_j)/2, \quad (3.8)$$

$$\mu_{1/2+\lambda} - \mu_{1/2} < (a_{j+1}-a_j)/2. \quad (3.9)$$

Note that  $\lambda$  exists since the median of  $F$  is unique and the set  $\{p: 0 < p < 1, \text{ the } p\text{-quantile is not unique}\}$  is countable.

Let  $x_{i,k_j}^*$  and  $y_{i,k_j}^*$  denote the  $i$ -th order statistics in  $B_{k_j, m}^{\circ}$  and  $B_{k_j+m, m}^{\circ}$  respectively. By lemma 3, for all  $k_j^{\circ}$ ,  $1 \leq j \leq q$ ,

$$x_{[m/2+m], k_j^{\circ}}^* \rightarrow \mu_{(1/2)+\lambda} + a_j, \text{ a.s.} \quad (3.10)$$

$$y_{[m/2-m], k_j^{\circ}}^* \rightarrow \mu_{(1/2)-\lambda} + a_{j+1}, \text{ a.s.} \quad (3.11)$$

By (3.8) and (3.9), with probability one for  $n$  large,

$$x_{[m/2+m], k_j^{\circ}}^* < y_{[m/2-m], k_j^{\circ}}^*, \text{ a.s.} \quad (3.12)$$

Now consider the order statistics in the combined sample  $B_{k_j^{\circ}, m}^{\circ} \cup B_{k_j^{\circ}+m, m}^{\circ} = A_{k_j^{\circ}, m}^{\circ}$ . By (3.12), and the fact that  $[m/2 + \lambda m] + [m/2 - \lambda m] \leq m$ , at least  $[m/2 + \lambda m]$   $x_i$ 's can be found in  $B_{k_j^{\circ}, m}^{\circ}$  which are the first  $m$  order statistics of  $A_{k_j^{\circ}, m}^{\circ}$ . Therefore, with probability one for  $n$  large, we have

$$S_{k_j^{\circ}, m}^{\circ} \geq [m/2 + \lambda m] - [m/2 - \lambda m] \geq \lambda m. \quad (3.13)$$

The lemma is proved.

#### 4. PROOF OF THEOREM AND SOME APPLICATIONS.

Proof of Theorem.

Let  $k_j^0$  be defined by (2.3). For fixed  $j$ , by lemma 4, we have, with probability one for  $n$  large,

$$S_{k_j^0, m}^2 / m \cong m \lambda^2.$$

By the definition of  $D_n$ ,  $k_j^0 \in D_n$ ,  $j = 1, \dots, q$ . Further, by the definition of  $D_{j,n}$ 's and  $n \gg m$ , it follows that with probability one for large  $n$ ,  $k_1^0, \dots, k_q^0$  belong to different  $D_{j,n}$ 's, which in turn implies

$$\hat{q} \cong q. \quad (4.1)$$

with probability one for large  $n$ . On the other hand, from Lemma 2 and  $C_n \gg \log n$ , it follows that with probability one for  $n$  large  $S_{k, m}^2 / m < C_n$ . Further take  $0 < \varepsilon < 1/2$  and write

$$K_n = \{k : m \leq k \leq n - m, |k/n - t_j| \geq (m/n)(1 + \varepsilon) \text{ for } j = 1, \dots, q\}. \quad (4.2)$$

Then, since  $n \gg m_n$ , and  $A_{k, m}$  consists of iid. random variables, it follows that with probability one for  $n$  large

$$D_n \cap K_n = \phi. \quad (4.3)$$

By the definition of  $D_{j,n}$  and the fact that  $2m(1 + \varepsilon) < 3m$ , we have with probability one for large  $n$

$$\hat{q} \leq q. \quad (4.4)$$

Combining (4.1) to (4.4), it follows that  $\hat{q} = q$  with probability one for  $n$  large, and

$$\lim_{n \rightarrow \infty} \hat{t}_j = t_j \quad j = 1, 2, \dots, q.$$

Thus, the theorem is established.

Our procedure can also be applied to the detection of change points of scale parameters. What we need to do is merely to consider  $\log |x(t)|$  instead of  $x(t)$ .

Likewise this method can be applied to the detection of change points of location parameters in a circle, i.e. directional data. Lombard (1986) proposed this problem and discussed the related testing and estimating problem when the number of change points is known. Our technique can handle the case in which this number is unknown.

For directional data, fix an origin and an axis and let  $\{x(t), 0 \leq t \leq 1\}$  be the angle process, where  $0 \leq x(t) < 2\pi$ . Assume that  $0 = t_0 < t_1 < \dots < t_q < t_{q+1} = 1$ ,  $t_1, \dots, t_q$  are the change points, so that  $x(t) \sim F(x - a_j)$  if  $t_{j-1} < t < t_j$ ,  $j = 1, \dots, q+1$ , where  $a_j \not\equiv a_{j+1} \pmod{2\pi}$ ,  $j = 1, \dots, q$ ,  $a_1, \dots, a_{q+1}$  denote the angular location parameters. The speciality of the angle process is that the angle is always in  $[0, 2\pi)$ . So the angle of  $\alpha + \beta$  equals  $\alpha + \beta - 2\pi$  if  $0 \leq \alpha, \beta < 2\pi$  and  $\alpha + \beta \geq 2\pi$ . A closer look at our proof reveals that the speciality has an effect only on (3.10) and (3.11) when  $\mu_{(1/2+\lambda)} + a_j = 0 \pmod{2\pi}$  or  $\mu_{(1/2-\lambda)} + a_{j+1} = 0 \pmod{2\pi}$  for some  $j$ ,  $1 \leq j \leq q$ . Therefore, if continuous distribution  $F(x)$  has a unique median and  $\mu_{1/2} + a_j \not\equiv 0 \pmod{2\pi}$  for all  $j$ ,  $1 \leq j \leq q$ , our procedure can be used to detect and estimate the change points of directional data. In practical applications, the origin and axis are chosen after gathering the data, in a way to make the analysis more convenient.

#### BIBLIOGRAPHY

- [1] Csörgö, M. and Horváth, L. (1988), Nonparametric methods for change point problems. *Multivariate Analysis* VII. to appear.
- [2] Hoeffding, W. (1963). Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, Vol.58, 13-30.
- [3] Krishnaiah, P. R. and Miao, B. Q. (1988). Review about estimates of change point. *Multivariate Analysis* VII. to appear.
- [4] Lombard, F. (1986). The change-point problem for angular data: a nonparametric approach. *Technometrics*, Vol. 28, No. 4, 391-397.
- [5] Pettitt, A. N. (1980). A simple cumulative sum type statistics for the change-point problem with zero-one observation. *Biometrika* Vol.67, 79-84.
- [6] Zacks, S. (1983). Survey of classical and Bayesian approaches to the change-point problem: Fixed sample and sequential procedures of testing and estimation. Recent advances in statistics, ed. M. H. Rizvi, New York: Academic Press, 245-269.